# A collective agenda on AI for Earth sciences

Manil Maskey, Ph.D.

AI for Good
February 16, 2022

# What I hope to do today

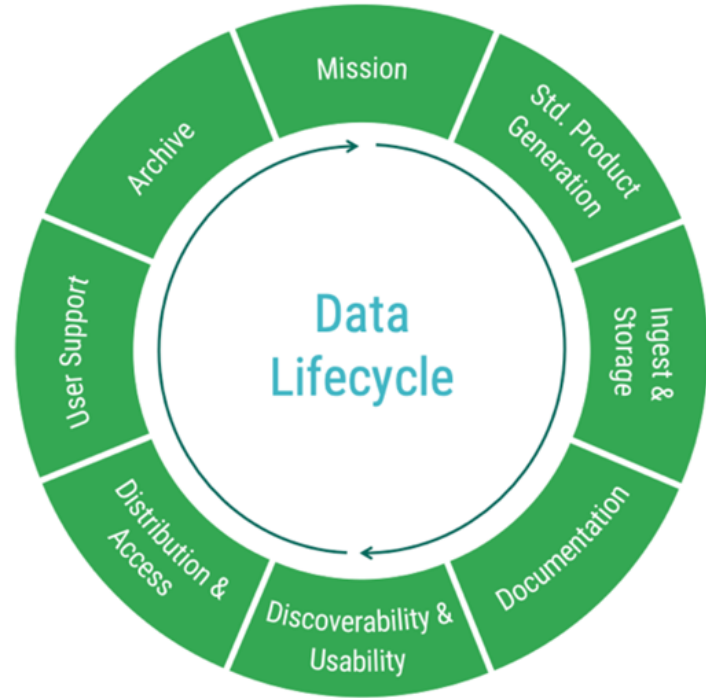Data systems perspective on AI for Earth science

- AI for core data systems services
    - Search
    - Knowledge discovery
- Enabling AI to advance Earth science
    - Data (labeled training data) is the proprietary differentiator
    - Transitioning AI models to Production
    - Citizen science

# NASA's Earth Science Data Systems Program

Single largest repository of Earth Science Data

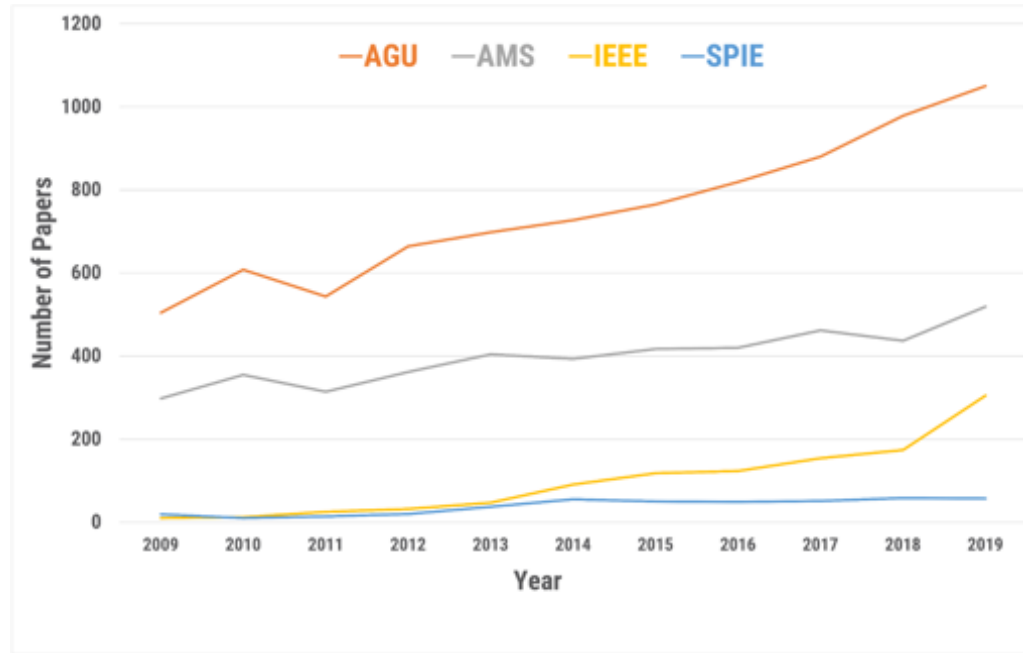Manages NASA's Earth science data through the entire data life cycle

# AGU topic trend
## Number of journal articles per keyword by year



2009

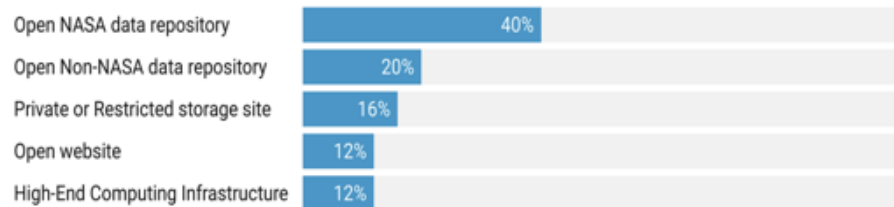# Publication trend - ML in Earth sciences



Rapid adoption of AI/ML by Earth science researchers
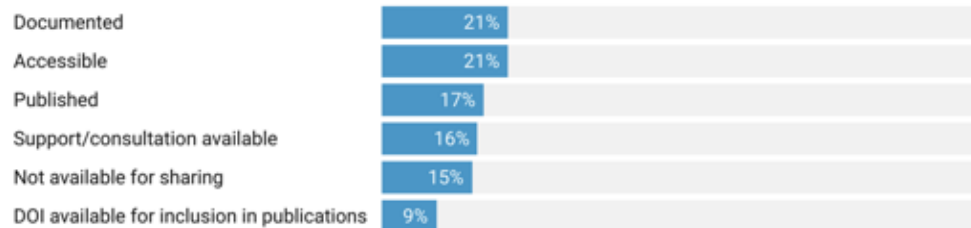
Virts et al. (2020)

Maskey et al. (2020)
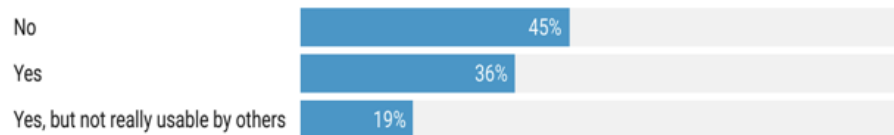
# NASA Science survey - AI and data
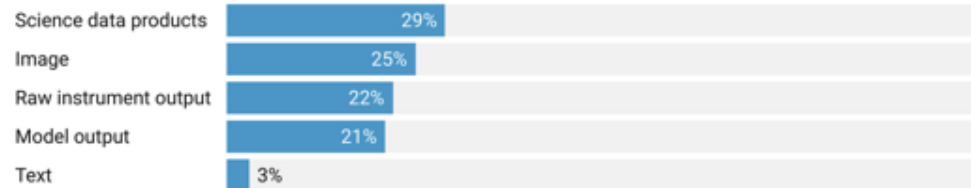
## Source of data used

| | |
|---|---|
| Open NASA data repository | 40% |
| Open Non-NASA data repository | 20% |
| Private or Restricted storage site | 16% |
| Open website | 12% |
| High-End Computing Infrastructure | 12% |

## How re-usable is your training data?

| | |
|---|---|
| Documented | 21% |
| Accessible | 21% |
| Published | 17% |
| Support/consultation available | 16% |
| Not available for sharing | 15% |
| DOI available for inclusion in publications | 9% |

## Is there a catalog of training data for your use?

| | |
|---|---|
| No | 45% |
| Yes | 36% |
| Yes, but not really usable by others | 19% |

## What type of data do you use for AI?

| | |
|---|---|
| Science data products | 29% |
| Image | 25% |
| Raw instrument output | 22% |
| Model output | 21% |
| Text | 3% |

## How did you construct training data?

| | |
|---|---|
| I created my own training data | 61% |
| I am re-using existing training data | 22% |
| I modified existing training data | 17% |

## Amount of effort required to prepare data for AI?

| | |
|---|---|
| Very low | 8% |
| Low | 13% |
| Average | 25% |
| High | 25% |
| Very high | 29% |

# NASA Science survey - AI and data

## Source of data used

| | |
|---|---|
| Open NASA data repository | 40% |
| **Open Non-NASA data repository** | 20% |
| Private or Restricted storage site | 16% |
| **Open website** | 12% |
| **High-End Computing Infrastructure** | 12% |

## How re-usable is your training data?

| | |
|---|---|
| Documented | 21% |
| Accessible | 21% |
| Published | 17% |
| Support/consultation available | 16% |
| Not available for sharing | 15% |
| DOI available for inclusion in publications | 9% |

## Is there a catalog of training data for your use?

| | |
|---|---|
| **No** | 45% |
| Yes | 36% |
| **Yes: but not really usable by others** | 19% |

## What type of data do you use for AI?

| | |
|---|---|
| **Science data products** | 29% |
| **Image** | 25% |
| Raw instrument output | 22% |
| Model output | 21% |
| Text | 3% |

## How did you construct training data?

| | |
|---|---|
| **I created my own training data** | 61% |
| I am re-using existing training data | 22% |
| I modified existing training data | 17% |

## Amount of effort required to prepare data for AI?

| | |
|---|---|
| Very low | 8% |
| Low | 13% |
| Average | 25% |
| **High** | 25% |
| **Very high** | 29% |

# Maximizing Knowledge Discovery

# Why?

Increasing Earth science data archives require non-traditional approaches to data management

Data driven technologies (AI) to provide advanced search capabilities

Machine learning-based approach - provide automated detection of Earth science events from image archives

Catalog of events can provide a novel way to explore large archives of data

Discover and explore Earth science data archives around events using machine learning (ML) techniques

Search

Browse Collections

Features

- Map Imagery
- Near Real Time
- Customize

Keywords

Platforms

Instruments

Organizations

Projects

Processing levels

fire

13 Matching Collections

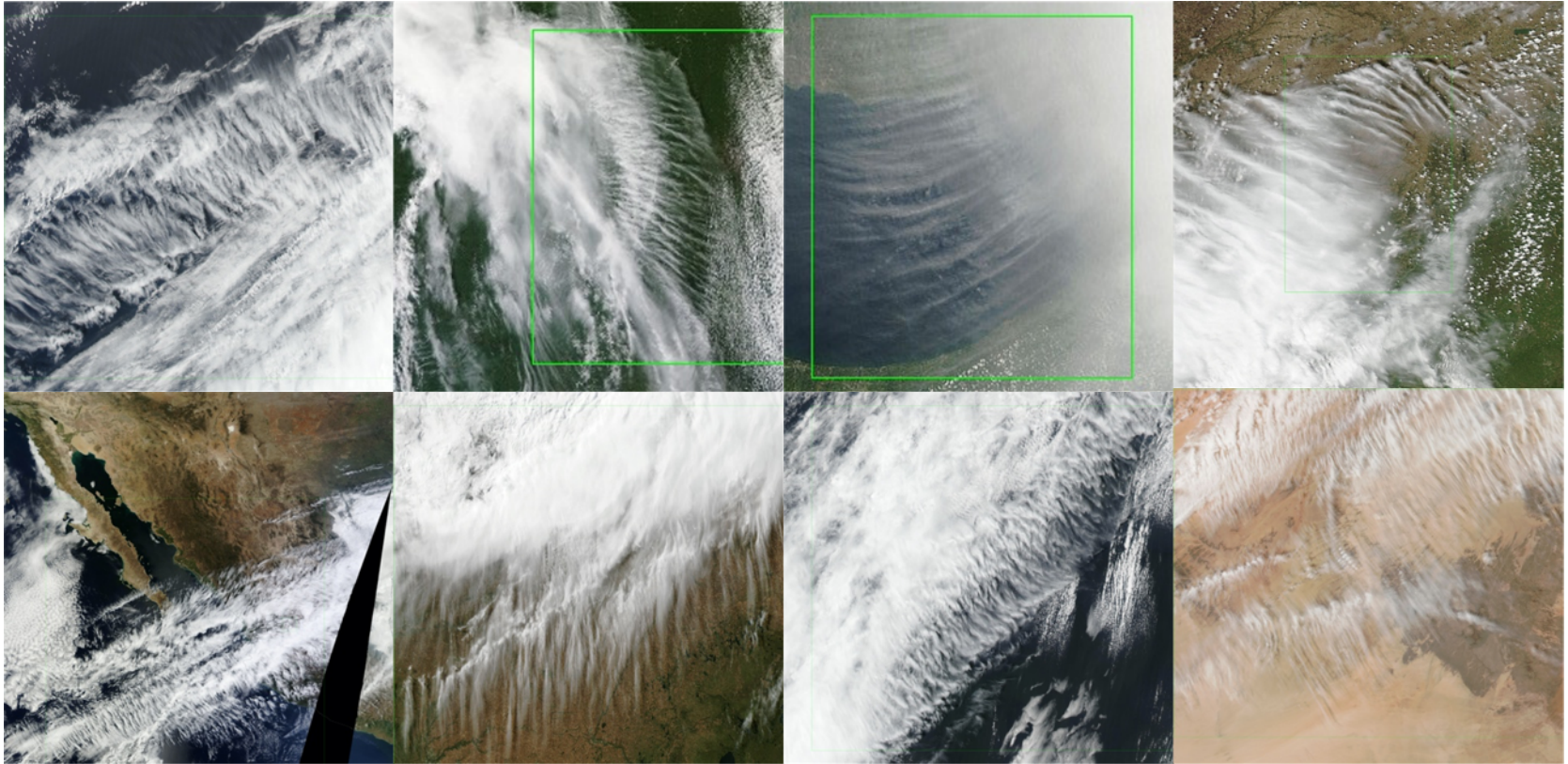Sort by: Relevance    Only include collections with granules    Include non-EOSDIS collections

High latitude dust

# Transverse cirrus bands

—

Welcome to the

# Phenomena Detection Portal

We are using machine learning for real-time detection of Earth science phenomena:

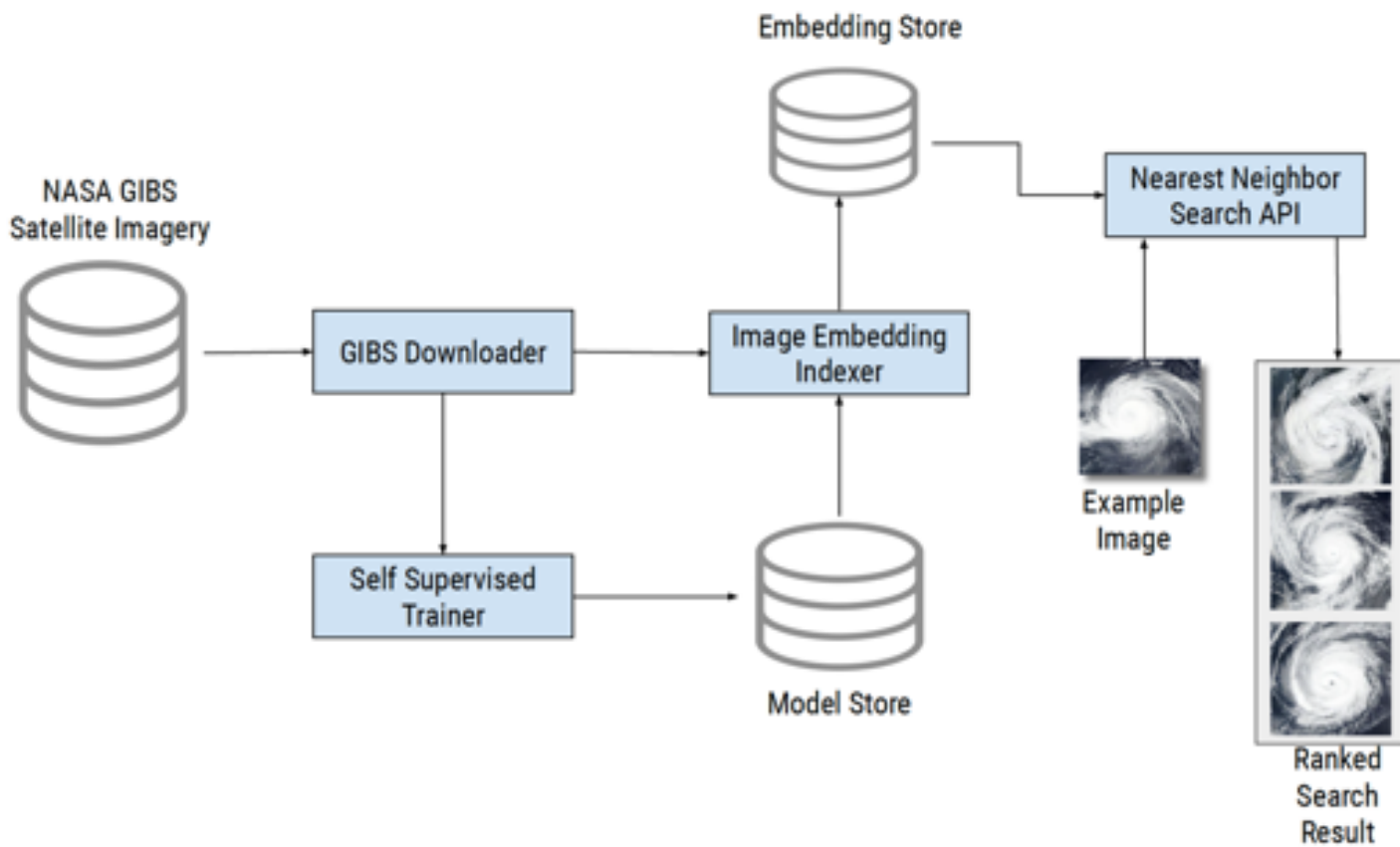| Types | Detections | Confidence score |
|---|---|---|
| **03** | **98,627** | **89.61%** |
| so far | and counting | on average |

**Start exploring**     Learn more

# Search by Example

Scaled SSL

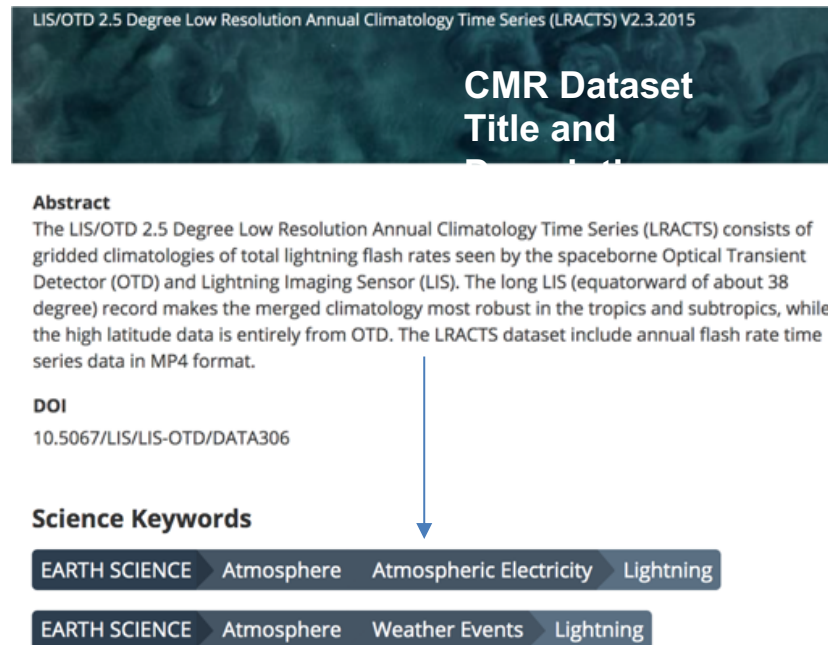# Augment data stewardship processes

# Why?

Assigning science keywords is currently a manual process, which is prone to human error and inconsistencies.

Metadata managed across a network of multiple data centers (i.e. keywords not assigned by a central entity)

Keywords may be assigned by non-subject matter experts (SMEs)

Improve metadata quality

Provide objective and consistent approach to keyword assignment

LIS/OTD 2.5 Degree Low Resolution Annual Climatology Time Series (LRACTS) V2.3.2015

**CMR Dataset Title and**

**Abstract**

The LIS/OTD 2.5 Degree Low Resolution Annual Climatology Time Series (LRACTS) consists of gridded climatologies of total lightning flash rates seen by the spaceborne Optical Transient Detector (OTD) and Lightning Imaging Sensor (LIS). The long LIS (equatorward of about 38 degree) record makes the merged climatology most robust in the tropics and subtropics, while the high latitude data is entirely from OTD. The LRACTS dataset include annual flash rate time series data in MP4 format.

**DOI**
10.5067/LIS/LIS-OTD/DATA306

**Science Keywords**

| EARTH SCIENCE | Atmosphere | Atmospheric Electricity | Lightning |

| EARTH SCIENCE | Atmosphere | Weather Events | Lightning |

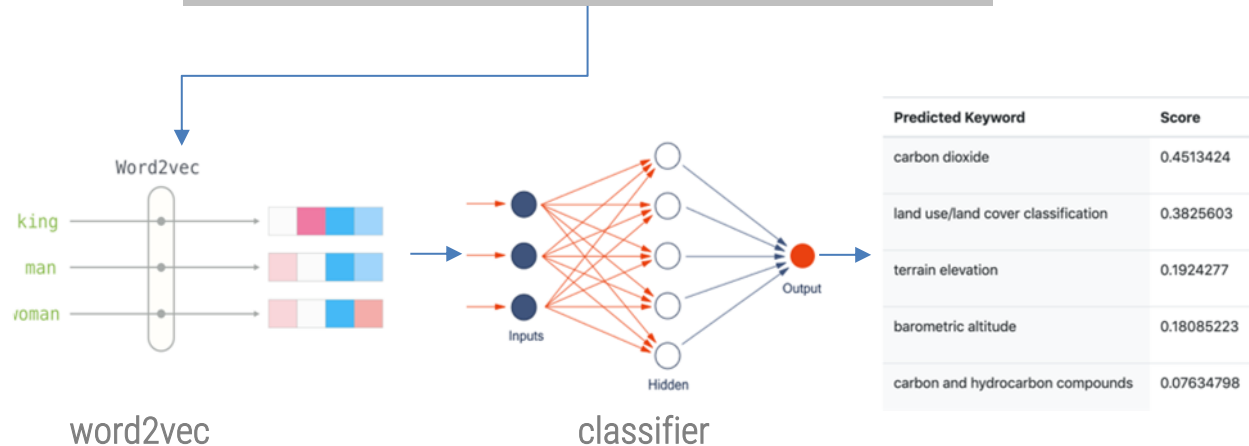# Approach – build word embeddings
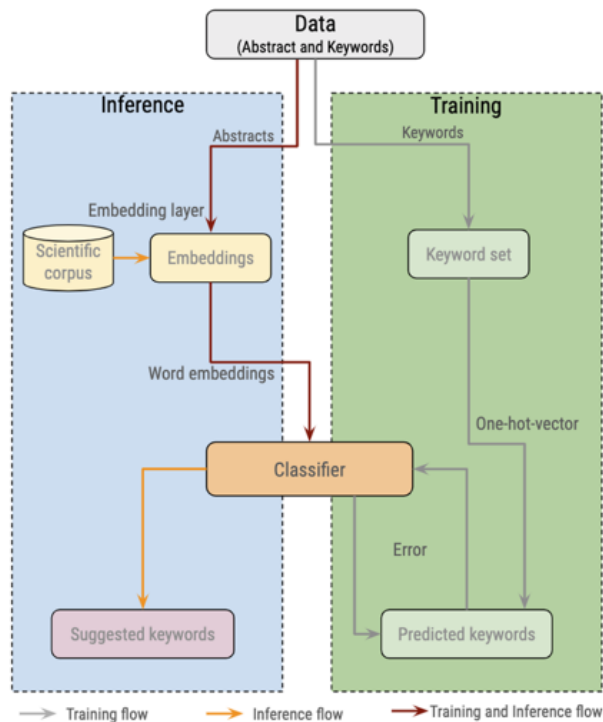


**88,410** documents

**530 million** words

**5.5 million** unique words

# Automated keyword assignment



Version 7.3 is the current version of the data set. Version 3.5 is no longer available and has been superseded by Version 7.3. This data set is currently provided by the OCO (Orbiting Carbon Observatory) Project. In expectation of the OCO-2 launch, the algorithm was developed by the Atmospheric CO2 Observations from Space (ACOS) Task as a preparatory project, using GOSAT TANSO-FTS spectra. After the OCO-2 launch, "ACOS" data are still produced and improved, using approaches applied to the OCO-2 spectra. The "ACOS" data set contains Carbon Dioxide (CO2) column averaged dry air mole fraction for all soundings for which retrieval was attempted. These are the highest-level products made available by the OCO Project, using TANSO-FTS spectral radiances, and algorithm build version 7.3. The GOSAT team at JAXA produces GOSAT TANSO-FTS Level 1B
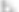
word2vec

classifier

| Predicted Keyword | Score |
|---|---|
| carbon dioxide | 0.4513424 |
| land use/land cover classification | 0.3825603 |
| terrain elevation | 0.1924277 |
| barometric altitude | 0.18085223 |
| carbon and hydrocarbon compounds | 0.07634798 |

# Keyword recommender

# Hurricane intensity estimation system

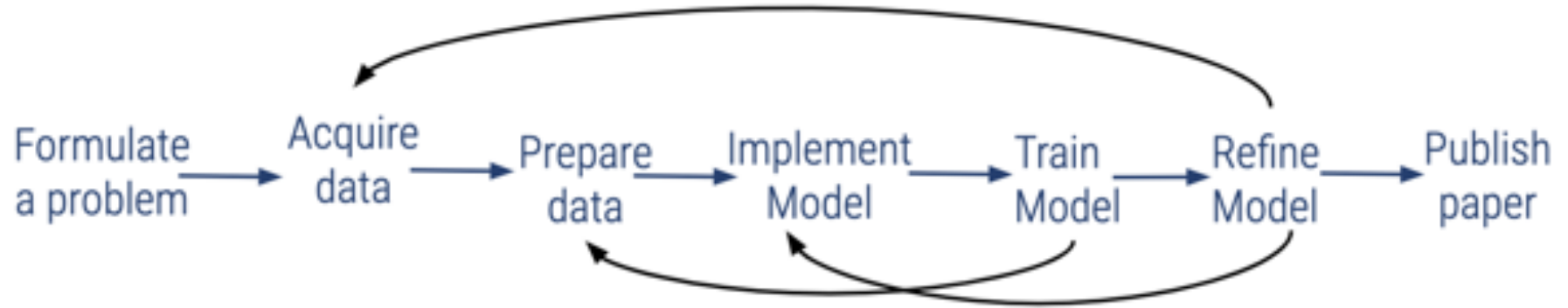# AI and satellite imagery to estimate hurricane wind speed

Hurricane Earl, 2010



Adapted from Stevenson et al. (2014). Time series of satellite-derived intensity estimates (circles) for Hurricane Earl (2010), added to best track intensities and lightning flash rate time series.

# ML in literature



Formulate a problem → Acquire data → Prepare data → Implement Model → Train Model → Refine Model → Publish paper

# We have a model…..now what?

Going extra mile

Interpretability + model inspection

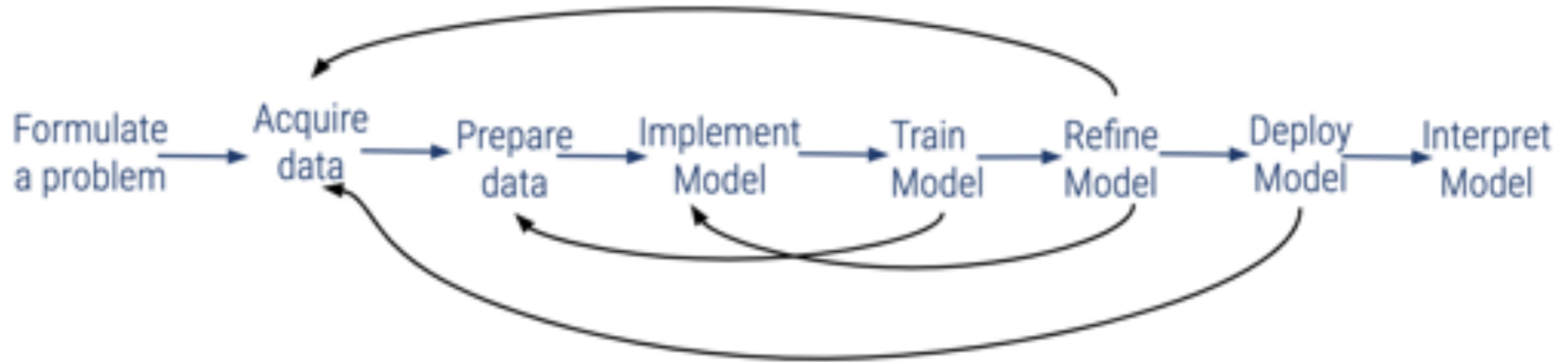Interpret prediction data – prediction output maybe just numbers

Questions:

Does the model confidence remain the same over time?
How do you maintain?
How do you complete the loop with new training data?

# ML lifecycle - iterative



Formulate a problem → Acquire data → Prepare data → Implement Model → Train Model → Refine Model → Deploy Model → Interpret Model

# Deployment to production

## Performance requirements

Metrics and baselines with initial models

Monitor over time

## Back-testing

Model and software will change

Testing model changes on historical data

Run current production model to baseline performance
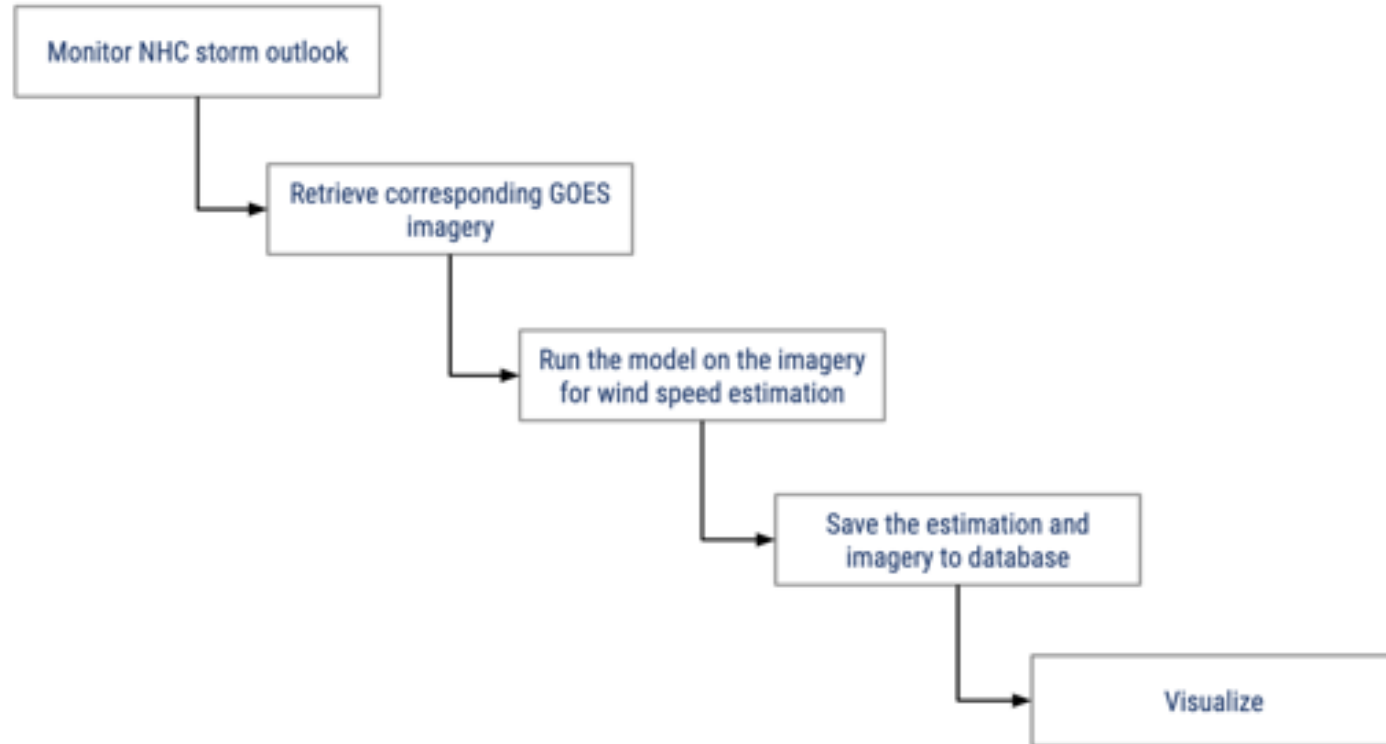
Run new models, competing for production

## Now-testing

Testing of production model on latest data

Can we get early warning that the model may be faltering?

- Content drift: training data exploited by model are subtly changing with time

# Workflow

Deep Learning-based Hurricane
Intensity Estimator

Applying machine learning to objectively estimate
tropical cyclone intensity.

Explore   or   Read more        Don't show again

Using community to advance model development

# Data science competitions

- Benchmark datasets and challenge problems have played an important role in driving progress in AI

  - Enables rigorous performance comparison

- Foster the learning of best practices

- Stimulate the abilities in problem-solving

- Encourage creativity and group work

- Give learners the chance to interact with new platforms and algorithms

- Citizen science

# Data science competition

**"Wind-dependent Variables: Predict Wind Speeds of Tropical Storms"**

- Leverage community to enhance solution to existing problem using open data
- Test whether high-quality datasets produces better models via open competition


- **Industry** partnership

- **733** participants

- **2756** entries

# Data science competition

## The Results

Over 700 participants stepped up to this important challenge, generating more than 2,700 entries. **Each of the top three models achieved at least a 50% reduction in Root Mean Square Error (lower is better) as compared to the existing model!**

### Root Mean Square Error (RMSE)



*Source: Maskey, et al. (2020)

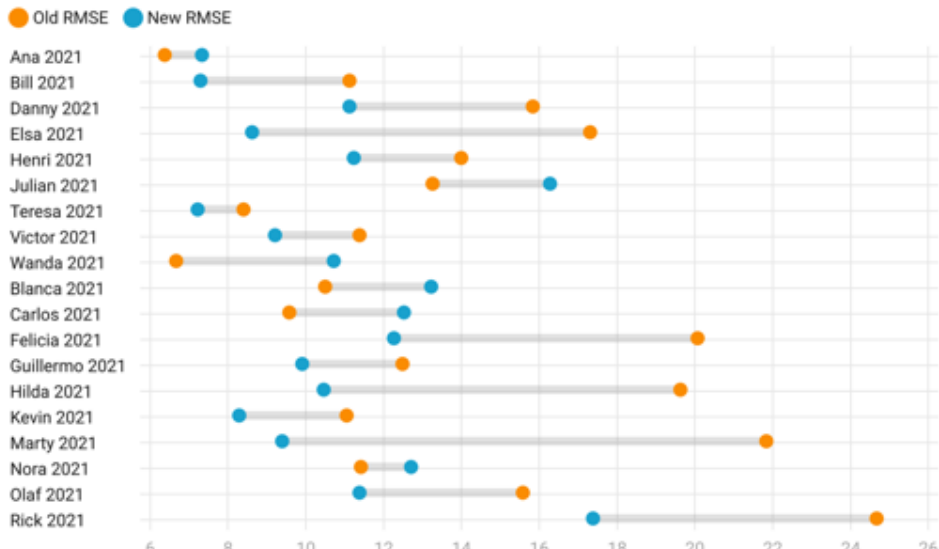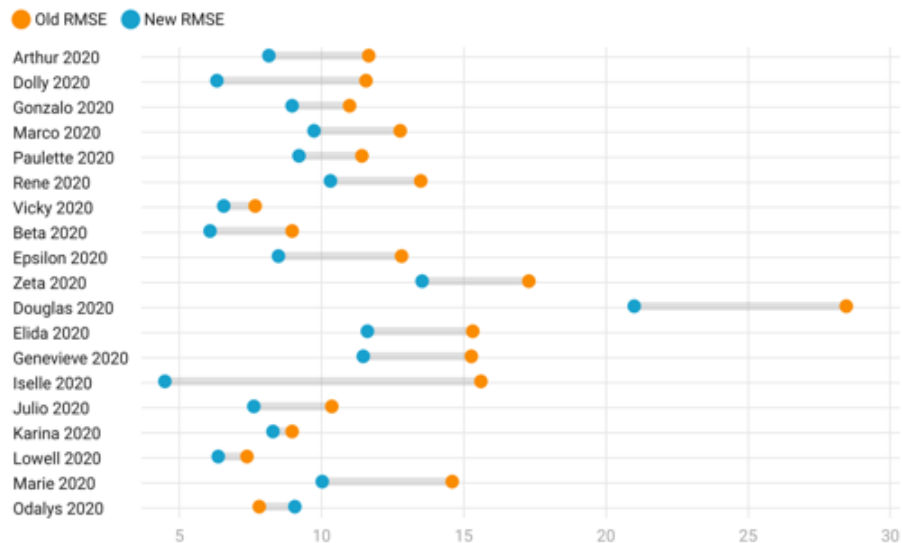| | First Place | Second Place | Third Place | Benchmark | Hurricane Intensity Estimator* |
|---|---|---|---|---|---|
| RMSE | 6.26 | 6.42 | 6.46 | 11.82 | 13.62 |

Winning solutions were able to take advantage of the relative timing of images in a storm sequence to produce targeted wind speed estimates based on temporal trends. As a result, these solutions can help to improve disaster readiness and response efforts around the world by equipping response teams with more accurate and timely wind speed measurements. All of the prize-winning solutions from this competition are linked below and made available for anyone to use and learn from.

# New Model – in depth analysis

# New Model – in depth analysis
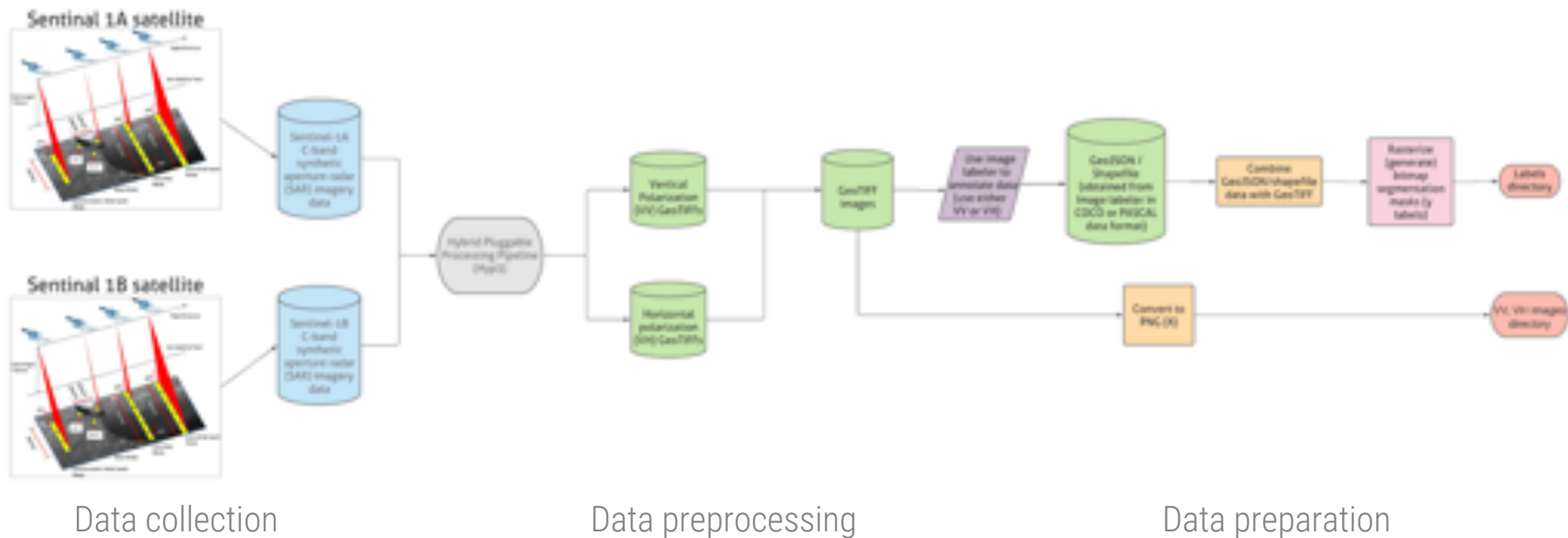
# Flood extent detection

# Flood extent detection

- A major natural disaster

- Widespread damage – property, agriculture

- Displacement, insurance, long-term socio-economic consequences

- Causes:
  - Persistent rainfall
  - Severe storm
  - High-tides
  - Storm surge from cyclones

# Problem

- Detecting flood extent is difficult

- Monitoring extent of flood events in-situ – hazardous to operate in a disaster zone

- Potential solution:
    - Remote sensing in conjunction with ML has been used in the community to monitor these events

- Need:
    - Large amounts of clean and labeled data

# Data acquisition



Data collection                    Data preprocessing                    Data preparation

# Data labeling

- 6 Atmospheric science/Earth science students

- 2 Domain scientists

- Training sessions

- Validation

**https://impact.earthdata.nasa.gov/labeler/**

# Labeled data

~66k images (33k VV + 33k VH images including swath gap artifacts)

Native resolution : 5x20m

## Train (24300):
- Nebraska (1741 sq. km.) (~43 %)
- North Alabama (13789 sq. km.) (~43%)
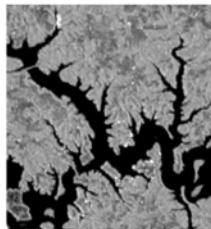- Bangladesh (7150 sq. km.) (~13%)

## Validation (6500):
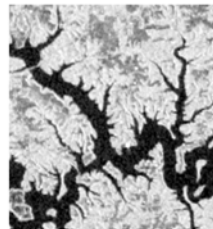- Florence (7197 sq. km.)

## Test (1600):
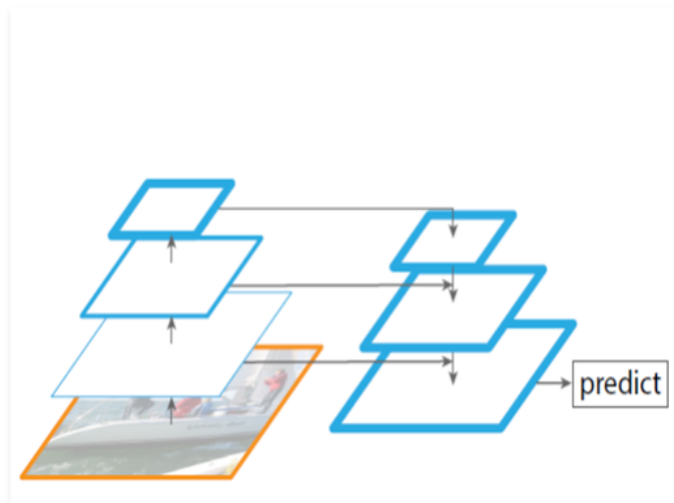- Red River North (6746 sq. km.)
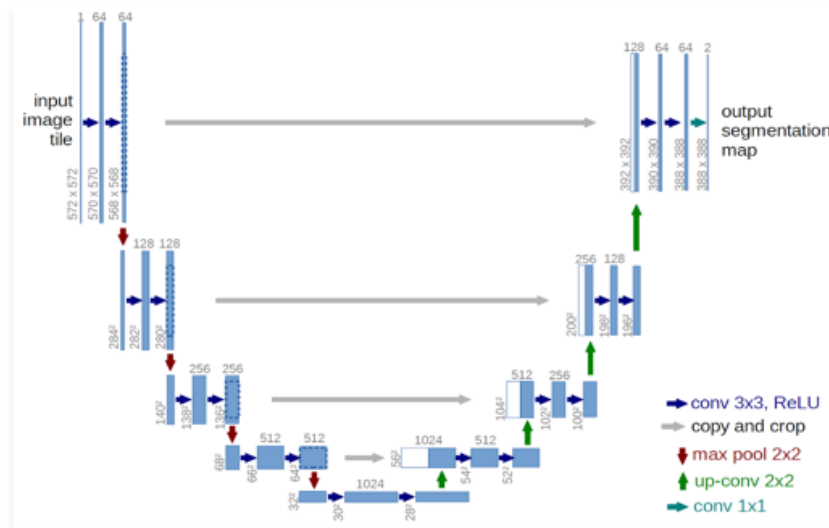


VV                   VH

# Benchmark model development

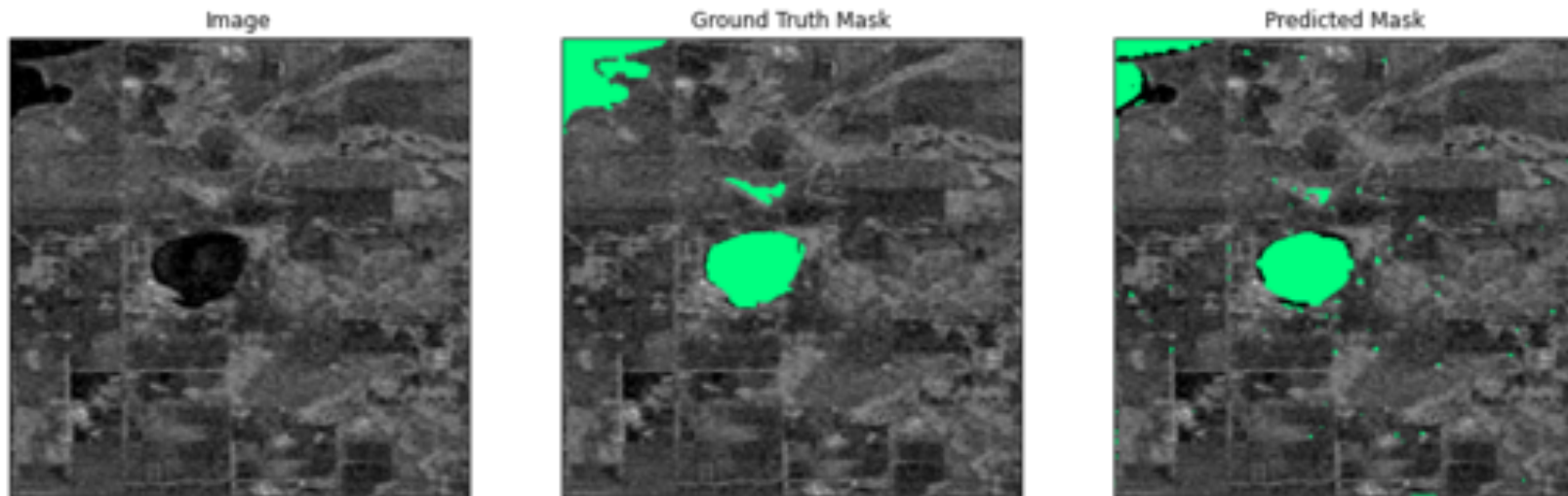Baseline models: Vanilla FPN and U-Net (ResNet 50 encoder). random selection, 4:1 Train-Validation split



Feature Pyramid Network (Resnet 50 encoder)



U-Net (Resnet 50 encoder)

# Benchmark model development

Visual Results



Sample segmentation

# Leveraging citizen science for optimal model

Finding optimal modal is an exhaustive task

ML Competition in collaboration with IEEE

- 137 participants
- More than 200 submissions
- Codalab platform

- **Phase 1 (Development):** Participants are provided with training data (which includes reference data) and validation data (without reference data until phase 1 concludes) to train and validate their algorithms. Participants can submit prediction results for the validation set to the codalab competition website to get feedback on the performance from April 15 to May 14, 2021. The performance of the best submission from each account will be displayed on the leaderboard.
- **Phase 2 (Test):** Participants receive the validation set reference data for model tuning and test data set (without the corresponding reference data) to generate predictions and submit their binary classification maps in numpy array format from May 15 to June 30, 2021. After evaluation of the results, three winners will be announced on July 1, 2021.

# Winning solutions



Team Arren, IOU: 0.7681

# Summary

AI enhanced enterprise data systems

- AI approaches that can efficiently operate as a part of the core of large-scale systems

- Before AI can be widely used in critical enterprise data systems, we need new robust pipelines to systematically manage AI lifecycle

- A flexible architecture that allows software systems and AI algorithms to evolve to take advantage of emerging trends in hardware and software and rapid model deployment

Thank you.